# Pose-Independent Automatic Target Detection and Recognition Using 3D Laser Radar Imagery

Alexandru N. Vasile and Richard M. Marino

■ **Although a number of object-recognition techniques have been developed to process terrain scenes scanned by laser radar (ladar), these techniques have had limited success in target discrimination, in part due to low-resolution data and limits in available computation power. We present a pose-independent automatic target detection and recognition system that uses data from an airborne three-dimensional imaging ladar sensor. The automatic target recognition system uses geometric shape and size signatures from target models to detect and recognize targets under heavy canopy and camouflage cover in extended terrain scenes. The system performance was demonstrated on five measured scenes with targets both out in the open and under heavy canopy cover, where the target occupied between 1% to 10% of the scene by volume. The automatic target recognition section of the system was successfully demonstrated for twelve measured data scenes with targets both out in the open and under heavy canopy and camouflage cover. Correct target identification was also demonstrated for targets with multiple movable parts in arbitrary orientations. The system achieved a high recognition rate along with a low false-alarm rate. Immediate benefits of the presented work are in the area of automatic target recognition of military ground vehicles, in which the vehicles of interest may include articulated components with variable position relative to the body, and may come in many possible configurations. Other application areas include human detection and recognition for homeland security, and registration of large or extended terrain scenes.**

THREE-DIMENSIONAL (3D) laser radar (ladar) sensors produce range images that provide explicit 3D information about a scene. Lincoln Laboratory has actively developed the laser and detector technologies that make it possible to build a high-resolution three-dimensional imaging ladar sensor with photon counting sensitivity [1]. In support of the Jigsaw program sponsored by the Defense Advanced Research Projects Agency (DARPA), the Laboratory has built a functional 3D ladar sensor system with a $32 \times 32$ array of avalanche photodiode (APD) detectors operating in Geiger mode. Recent field tests using this Jigsaw ladar sensor produced high-quality 3D imagery of targets behind obscurants for extremely low signal levels [1].

The primary purpose of a ladar sensor is to record the 3D spatial signature of a target so that the particular target can be identified. As an extension of the Jigsaw program, Lincoln Laboratory has developed a complete end-to-end automatic target detection and recognition (ATD/R) system. The implemented target detection and recognition algorithms use field data collected by the high-range-resolution Jigsaw ladar sensor, as well as some data sets taken with the previous GEN-III ladar sensor [2].

The primary goal of the ATD/R system is to accurately detect and recognize targets present in large terrain scenes, where the target may occupy less than 1% of the scene and have more than two hundred points on target. A secondary system goal was to demonstrate correct target identification with foliage occlusion greater than 70%. Another goal was to demonstrate correct identification of articulated targets, with multiple movable parts that are in arbitrary orientations. The above goals have to be met while achieving a high recognition rate (over 99%) along with a low false-alarm rate (less than 0.01%).

## Background on Target Detection

The problem of automatic target recognition in ladar range imagery has been an active topic of research for a number of years [3, 4]. Automatic target recognition (ATR) involves two main tasks: target detection and target recognition [5]. The purpose of target detection is to find regions of interest (ROI) where a target may be located. By locating ROIs, we can filter out a large amount of background clutter from the terrain scene, making object recognition feasible for large data sets. The ROIs are then passed to a recognition algorithm that identifies the target [5].

Target detection methods attempt to determine the presence of a target in a large data set by quickly filtering large portions of the scene prior to submitting the data to the recognition algorithm. In the ATR field, detection methods that can search a large data set and reduce it to a few ROIs are known as cueing algorithms [6]. The application of a cueing algorithm as a data-preprocessing step vastly reduces the time needed for target recognition.

Target detection approaches can be classified as image based and model based [7]. The traditional image-based approach uses template matching; the target is separated from its surrounding area by extracting a silhouette based on a target image template [8]. However, silhouette extraction algorithms do not reliably recover the true silhouette from real imagery, thus seriously degrading the robustness of target detection [8]. In general, the template approach suffers from the complexity in finding the silhouette in the image, as well as the complexity of creating the template database [7].

With significant improvements in ladar sensor resolution and increased computational power, detailed 3D structural information may be obtained from the data and used by model-based approaches. Traditional model-based approaches rely on boundary segmentation and planar surface extraction to describe the scene. Target detection is then performed through the use of trained neural networks or genetic algorithms [8–12]. One recent cueing algorithm that is applicable to large ladar data sets is the spin-image–based 3D cueing algorithm developed by O. Carmichael and M. Hebert [6].

Given an ROI, the recognition algorithm attempts to classify the particular target in a library of target models. The target models are used to represent a unique signature that is present in the target data set. There are numerous ways to encode the target models. For ladar data, where the scene data consists of an unstructured point cloud, object representation schemes can be divided into two categories: surface-based 3D model representations and shape-based two-dimensional (2D) model representations.

Surface-based 3D model representation schemes perform geometrical surface matching between a library of 3D surface models and a data scene. Traditional 3D geometrical feature-matching algorithms segment the target into simple geometric primitives, such as planes and cylinders, and record the relative spatial relationship of each geometric primitive in the target model [13–15]. The scene is then segmented in the same manner, and the library is searched for a group of primitive objects that have a spatial structure similar to the target model's [16, 17]. Recent methods have shown that planar patch segmentation is robust to noisy range data [18]. In addition, current 3D feature-grouping schemes have been proven to work even when the target is partially occluded [19].

An alternate approach to 3D geometric feature matching is to reduce the 3D recognition problem to a set of 2D recognition problems, in which the target signature is encoded by a shape-based 2D representation. The primary advantage of the shape-based recognition approach over 3D geometrical matching is that it can scale well to large data sets with high levels of clutter [3, 20]. In addition, the recognition algorithms can benefit from the tremendous amount of work done in the relatively mature field of 2D image analysis. Some recent algorithms that use shape-based representations are the
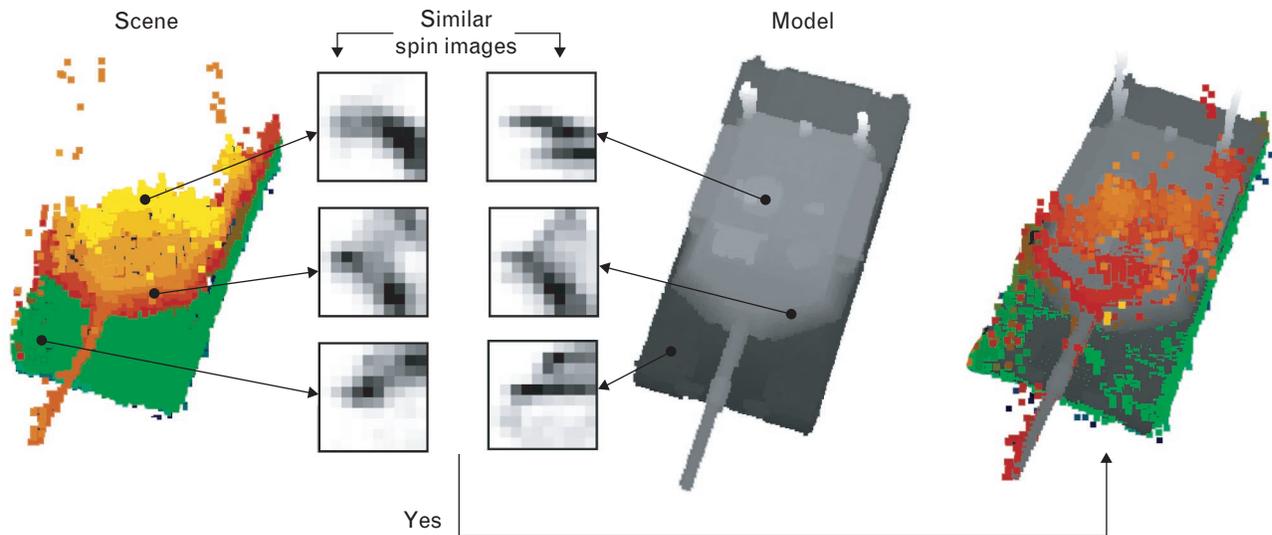
Scene    Similar
spin images    Model



**FIGURE 1.** Spin-image surface-matching concept. Given a target scene (with height color coded in green, red, and yellow), we can create a spin image for each scene point. Similarly, we can also create a spin image for each point in the model data set (with height color coded in shades of gray). For each scene spin image, we search through all the model spin images and find the best match. In this way, correspondences are found between the scene points and the model points. These correspondences can then be used to compute a three-dimensional (3D) transformation that aligns the scene and the model data sets.

contour-based algorithm of V. Shantaram et al. [21], the shape spectra algorithm of C. Dorai et al. [22], the surface signatures of S. Yamany et al. [23] and A. Johnson's spin-image algorithm [24].

After performing a literature review of the current techniques in target detection and recognition using ladar imagery, we found the spin-image–based detection and recognition algorithms to be most promising for processing our 3D ladar terrain data [25]. The remainder of this article is a discussion of the two main component areas—automatic target recognition and automatic target detection—of these spin-image–based algorithms.

**Automatic Target Recognition**

Given an ROI within a large-scale scene, the ATR algorithm attempts to identify a potential target from among the targets in a model library, or else it will report a none-of-the-above outcome. The recognition algorithm as well as the detection algorithm are based on Johnson's spin-image surface matching. We give here an overview of spin-image surface matching to provide a context for understanding the development of algorithms to follow.

*Spin-Image Surface Matching*

In the spin-image–based representation, surface shape is described by a collection of oriented 3D points with associated surface normals. Each 3D oriented point has an associated image that captures the global properties of the surface in an object-centered local coordinate system [24]. By matching images, we can determine correspondences between surface points, which results in surface matching. Figure 1 illustrates the spin-image surface-matching concept.

The image associated with each 3D oriented point is known as a spin image. A spin image is created by constructing a local coordinate system at an oriented point. By using this local coordinate system, we can encode the position of all the other points on the surface with two parameters: the signed distance in the direction of the surface normal and the radial distance from the surface normal. By mapping many of the surface points to this 2D parameter space, we can create a spin image at each oriented point. Since a spin image encodes the coordinates of the surface points with respect to a local coordinate system, it is invariant to rigid 3D transformations. Given that a 3D point can now be described by a
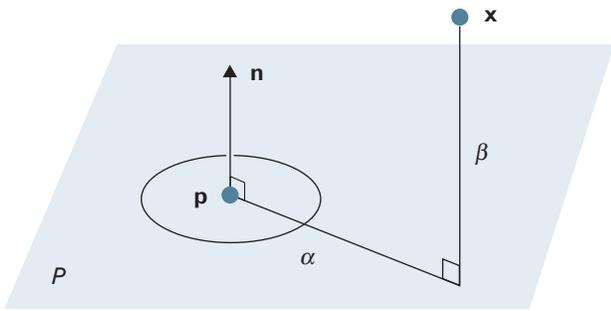
**FIGURE 2.** Constructing an oriented point basis for a 3D point **p**. Given an oriented point **p** in a tangent plane $P$ with unit normal **n**, a two-dimensional (2D) parameter space can be created that is invariant to pose. A point **x**, belonging to the same data set as point **p**, can be projected to this 2D parameter space. The quantity $\alpha$ is the distance from **x** to **p** perpendicular to the normal **n**, and $\beta$ is the signed distance from **x** to the plane $P$.

corresponding image, we can apply robust 2D template matching and pattern classification to solve the problem of surface matching and 3D object recognition [24].

The fundamental component for creating a spin image is the associated 3D oriented point. As shown in Figure 2, an oriented point defines a five-degree-of-freedom basis, using the tangent plane $P$ though point **p**, oriented perpendicular to the unit normal **n**.

Two coordinates can be calculated, given an oriented point: $\alpha$ is the perpendicular distance to the unit surface normal **n**, and $\beta$ is the signed perpendicular distance to the plane $P$ [24]. Given an oriented point basis O, we can define a mapping function $S_{O}$ that projects

3D points **x** to the 2D coordinates of a particular basis $(\mathbf{p}, \mathbf{n})$ as follows:

$$S_{O} : R^3 \rightarrow R^2$$
$$S_{O}(\mathbf{x}) \rightarrow (\alpha, \beta)$$
$$= \left( \sqrt{|\mathbf{x} - \mathbf{p}|^2 - [\mathbf{n} \cdot (\mathbf{x} - \mathbf{p})]^2}, \ \mathbf{n} \cdot (\mathbf{x} - \mathbf{p}) \right).$$

Applying the function $S_{O}(\mathbf{x})$ to all the oriented points in the 3D point cloud will result in a set of 2D points in $\alpha - \beta$ space. To reduce the effect of local variations in 3D point positions, we attach the set of 2D points to a 2D array representation grid. Figure 3 illustrates the procedure to create a 2D array representation of a spin image. To account for noise in the data, we linearly interpolate the contribution of a point to the four surrounding bins in the 2D array. By spreading the contribution of a point in the 2D array, bilinear interpolation helps to further reduce the effect of variations in 3D point position on the 2D array. This 2D array is considered to be the fully processed spin image.

The implemented surface-matching algorithm follows closely the procedure described in chapter 3 of Johnson's Ph.D. thesis [24]. The algorithm takes a scene data set along with a spin-image model library. The spin-image model library contains the ideal 3D ladar signatures of each target, derived from computer-aided design (CAD) models. Each 3D ladar model data set also has an associated spin-image database, with a corresponding spin image for each model 3D point.
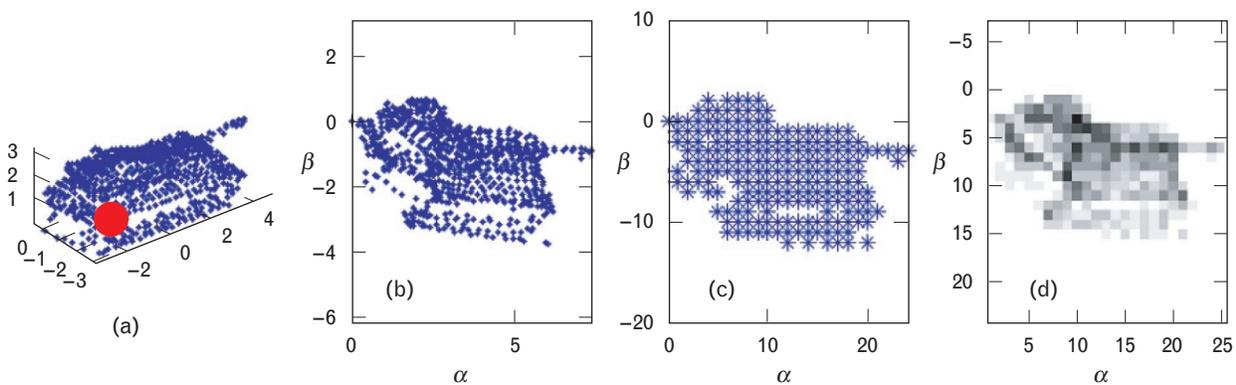


**FIGURE 3.** A 2D array representation of a spin image using bilinear interpolation. (a) Measurements of an M60 tank, in metric units. The red dot indicates the location of the 3D point used to create the example spin image. (b) Resulting mapping of the scene points in the $\alpha - \beta$ spin-map of the chosen 3D point, in metric units. (c) Spin image showing the non-zero bins after applying bilinear interpolation. (d) Spin image showing the bin values on a gray color scale. The darker bins indicate that a larger number of points were accumulated in those particular bins.

The spin-image model library is computed *a priori* to save online recognition time. The spin-image algorithm takes the scene data set and creates a spin-image database based on a subsampling of the points. The sampling ranges from 20% to 50% of all scene data points. The scene data points are not judiciously picked: the points are uniformly distributed across the given scene. Therefore, no feature extraction is performed to pick spin-image points.

The scene spin-image database is correlated to each model spin-image database within the model library. For a scene-to-model comparison, each scene spin image is correlated to all the model spin images, resulting in a distribution of similarity measure values. The correspondences obtained for each scene spin image to model spin-image database comparison are filtered by using a statistical data-based similarity-measure threshold. The above process is repeated for the rest of the scene spin images, resulting in a wide distribution of similarity measures.

Given the new distribution of similarity measures, a second similarity threshold is applied to remove un-

likely correspondences. The remaining correspondences are further filtered and then grouped by geometric consistency in order to compute plausible transformations that align the scene to the model data set. The initial scene-to-model alignment is refined by using a modified version of the iterative closest point (ICP) algorithm to obtain a more definite match. Figure 4 shows a detailed block diagram of the surface-matching process.

This particular surface-matching process is versatile, since no assumptions are made about the shape of the objects represented. Thus arbitrarily shaped surfaces can be matched without the need for initial transformations. This matching is particularly critical for our target recognition problem in which the target's position and pose within the scene are unknown. Furthermore, by matching multiple points between scene and model surfaces, the algorithm can eliminate incorrect matches due to clutter and occlusion.

The end result of spin-image–based surface matching is an optimal scene-to-model transformation, along with a recognition goodness of fit ($R_{\mathrm{GOF}}$) value between the scene and the model. The $R_{\mathrm{GOF}}$ of a comparison of
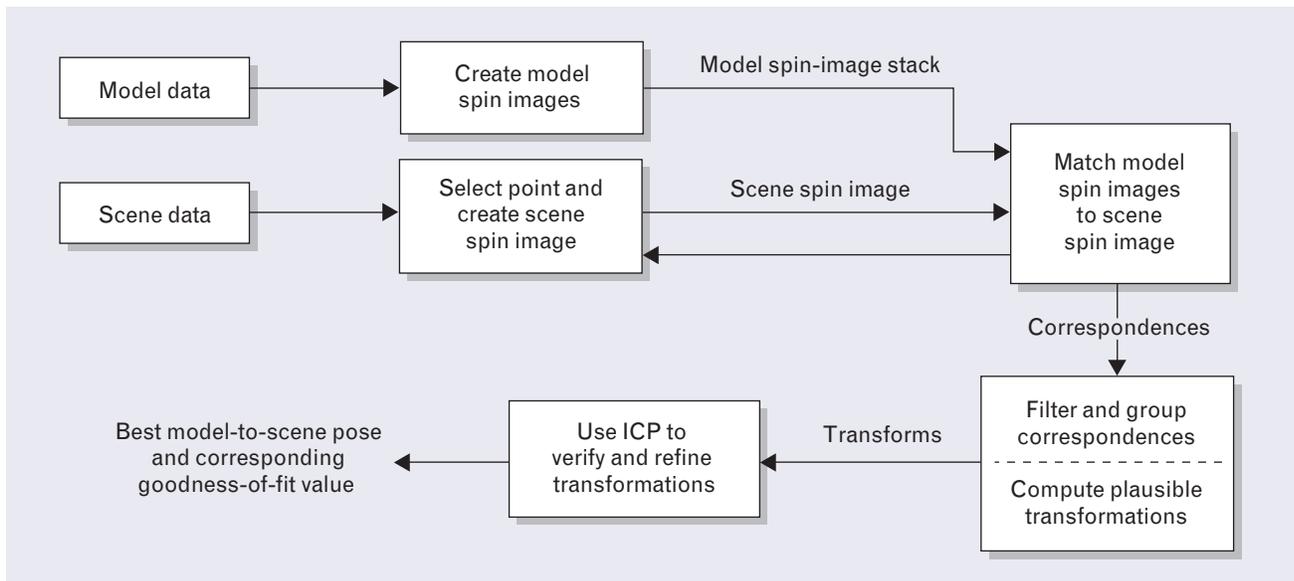


**FIGURE 4.** Surface-matching diagram. The matching process starts with a scene data set and a model data set. A spin-image stack is created for the model data set. For the scene data set, several points are randomly selected, and corresponding spin images are computed. For each scene spin image, we search through all the model spin images and find the best match. In this manner, correspondences are found between the scene and model points. After some filtering steps, the correspondences can then be used to compute an iterative closest point (ICP) 3D transformation that aligns the scene data set and the model data set. On the basis of this model-to-scene alignment, the process determines a goodness-of-fit value to score the scene-to-model match.

scene $s$ to model $m$ is defined as

$$R_{\mathrm{GOF}}(s,m) = \frac{(\theta^2 \cdot N_{pt})}{\mathrm{MSE}} , \qquad (1)$$

where $\theta$ is the fraction of overlap between the scene and the model as determined by the ICP algorithm, $N_{pt}$ is the number of plausible pose transformations found by the spin-image correlation process, and MSE is the mean-squared error as determined by the ICP algorithm. A higher $R_{\mathrm{GOF}}$ value indicates a higher level of confidence that the model matches the scene.

To quantify the recognition performance of a scene-to-model library comparison, we normalize the $R_{\mathrm{GOF}}$ to the sum of all the found $R_{\mathrm{GOF}}$ values. The normalized $R_{\mathrm{GOF}}$ that the scene $s$ correctly matches model $i$ in a model library *mlib* is defined as

$$\bar{R}_{\mathrm{GOF}}(s,mlib_i) = \frac{R_{\mathrm{GOF}}(s,mlib_i)}{\displaystyle\sum_{j=1}^{N} R_{\mathrm{GOF}}(s,mlib_j)} ,$$

where $N$ is the number of models in the model library *mlib*.

For each scene-to-model library comparison, the $\bar{R}_{\mathrm{GOF}}$ is split among the models and ranges from zero to one. For a given scene, the sum of the $\bar{R}_{\mathrm{GOF}}$ values over all the models in the model library adds up to one, unless a "none-of-the-above" outcome is reached. In the case of a "none-of-the-above" conclusion, the sum of the $\bar{R}_{\mathrm{GOF}}$ values equals zero, and each $\bar{R}_{\mathrm{GOF}}$ equals zero by definition.

The higher the value of $\bar{R}_{\mathrm{GOF}}$ for a scene-to-model comparison, the more likely it is that the model correctly matches the given scene. Thus the $\bar{R}_{\mathrm{GOF}}$ value that falls on each model represents a confidence measure that the model matches the scene.

## Results and Discussion

The ATR results presented here are divided into two sections. The main section is devoted to the non-articulated ATR results obtained from the comparison of twelve measured data scenes to a target model library consisting of ten target vehicles. A second, smaller section focuses on the results of a limited study of articulated ATR.

### Non-Articulated ATR Study

For the study of non-articulated ATR, we used the target model library that was developed under the Jigsaw program. The Jigsaw model library has approximately ten targets of interest, ranging from trucks and armored personnel carriers (APC) to tanks and missile launchers. Figure 5 shows the CAD models of the specific targets. The model library contains two large target classes, namely, APCs and tanks. The APC target class is composed of the BMP-1, BMP-2, BTR-70, and M2 vehicles. The tank class includes the M1A1, M60, and T72 tanks.

With the above CAD models, we constructed a target model library to simulate an ideal 3D ladar signature of each target. The simulated targets were then represented in the spin-image representation as 3D oriented points with associated spin images. We used the result-
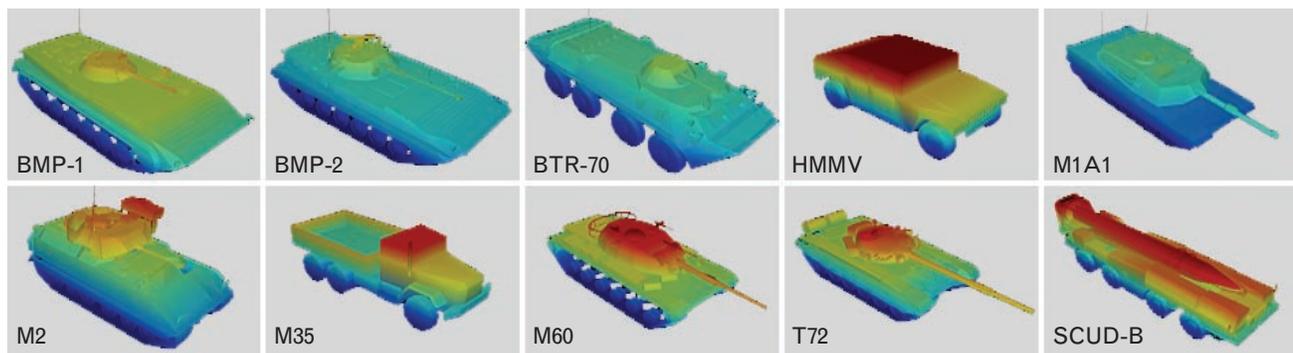


**FIGURE 5.** Examples of computer-aided design (CAD) target models, color coded by height, in the model library developed for the Jigsaw program. These models include trucks, armored personnel carriers, tanks, and missile launchers.

### Table 1. Resulting 3D Oriented Point Data Sets for the Given Target Models for Two Subsampling Voxel Sizes

| Subsampling voxel size (m) | Estimated surface resolution (m) | Number of Points in the Model Data Set | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | BMP-1 | BMP-2 | BTR-70 | HMMV | M1A1 | M2 | M35 | M60 | SCUD-B | T72 |
| 0.1 | 0.125 | 10,366 | 9692 | 11,444 | 5035 | 16,394 | 14,669 | 10,146 | 18,778 | 23,916 | 13,454 |
| 0.2 | 0.25 | 2239 | 2056 | 2453 | 1255 | 3761 | 3223 | 2368 | 4035 | 5213 | 2842 |

ing model spin-image library to compare the models to measured scenes in order to recognize and identify the scene target. Table 1 summarizes the resulting model data sets obtained from the 3D simulation for two voxel subsamplings.

Multiple scenes were analyzed to determine the recognition performance. A recognition confusion matrix was calculated as a measure of the recognition performance, showing the confidence measurement $\overline{R}_{\mathrm{GOF}}$ on the main diagonal and errors on the off diagonals [26]. Twelve scenes, each containing a target instance, were used to create the confusion matrix. Target truth was known prior to data collection. Measured data for the following targets were used: BMP-1, BTR-70, HMMV, M1A1, M2, M35, M60 and the T72. Figure 6 shows an orthographic projection of each of the twelve measured scene data sets.

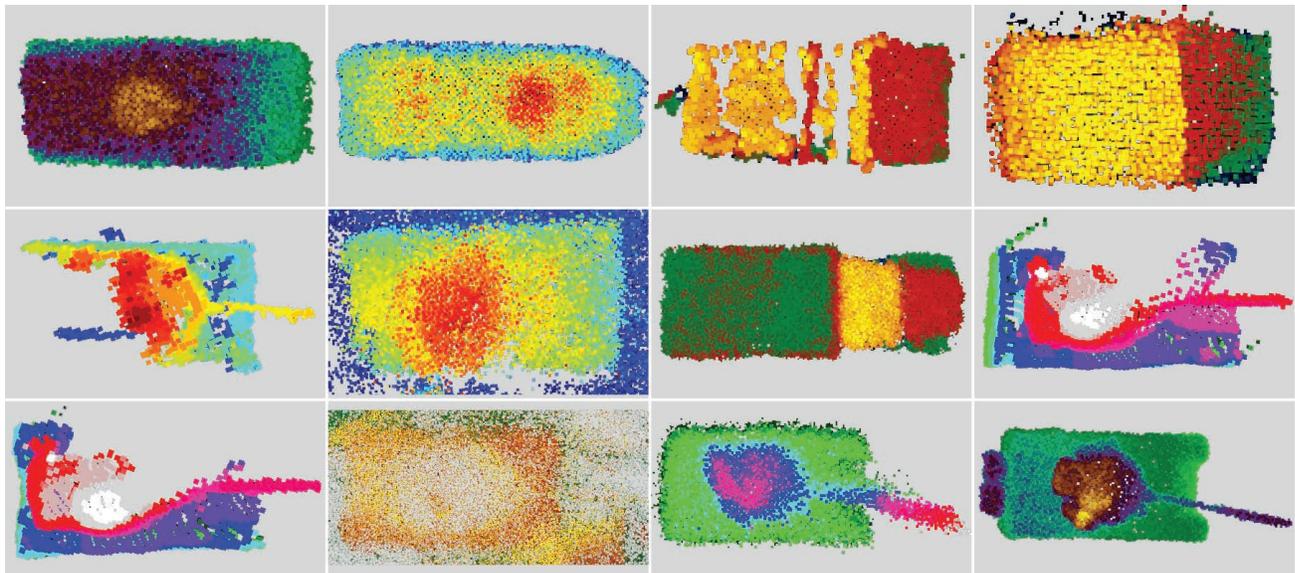Table 2 shows the recognition confusion matrix we



**FIGURE 6.** Orthographic view of the twelve measured scene data sets, color coded by height. These scenes served as target truth for comparisons with the model library.

### Table 2. Recognition Confusion Matrix*

| Field Data | Angular diversity | Angular views | BMP-1 | BMP-2 | BTR-70 | HMMV | M1A1 | M2 | M35 | M60-A3 | SCUD-B | T72 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | *Models* | | | | | | |
| **BMP-1** C5-F10-P03 | 10° | 16 | **0.61** | 0.38 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.01 |
| **BTR-70** C5-F10-P04 | 10° | 23 | 0.0 | 0.0 | **0.81** | 0.0 | 0.19 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **HMMV** RMF May 2002 | 15° | 4 | 0.0 | 0.0 | 0.0 | **1.0** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **HMMV** C8-F01-P10 | 30° | 10 | 0.0 | 0.01 | 0.0 | **0.92** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.07 |
| **M1A1** Eglin Dec 01 | 0° | 1 | 0.0 | 0.0 | 0.0 | 0.0 | **0.92** | 0.01 | 0.0 | 0.0 | 0.0 | 0.07 |
| **M2** C5-F13-P07 | 20° | 16 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | **1.0** | 0.0 | 0.0 | 0.0 | 0.0 |
| **M35** C5-F10-P05 | 15° | 12 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | **1.0** | 0.0 | 0.0 | 0.0 |
| **M60-A3** w/plow Huntsville May 2002 | 0° | 1 | 0.0 | 0.0 | 0.09 | 0.0 | 0.0 | 0.0 | 0.0 | **0.91** | 0.0 | 0.0 |
| **M60-A3** Huntsville May 2002 | 0° | 1 | 0.0 | 0.03 | 0.0 | 0.0 | 0.0 | 0.01 | 0.0 | **0.96** | 0.0 | 0.0 |
| **M60-A3** C05-F16-P10 | 10° | 12 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | **0.97** | 0.01 | 0.02 |
| **T72** C05-F00-P03 | 15° | 105 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | **1.00** |
| **T72** C20-F01-P03 | 15° | 29 | 0.0 | 0.0 | 0.0 | 0.0 | 0.13 | 0.0 | 0.0 | 0.0 | 0.0 | **0.87** |

\* Each row of the confusion matrix represents a scene-to-target model library comparison. Each cell in a row shows the resulting normalized $R_{GOF}$ that the target (with the identifying label shown in the top row) matches the scene (described at the beginning of the row). For each scene, the angular diversity and angular view are also shown in the first two columns to give a notional idea of the target coverage or obscuration.

obtained from the comparison of the model library to each of the twelve scenes. Each row of the confusion matrix represents a scene-to-model library comparison. For instance, the first row contains the comparison between a BMP-1 scene measurement and the model library. We see that the recognition confusion matrix resembles an identity matrix, which would be the ideal result. For all scene comparisons, the highest $\bar{R}_{GOF}$ value always falls on the target that matches the scene target truth. Furthermore, $\bar{R}_{GOF}$ has a value of zero for most of the remaining targets because the recognition algorithm found no match between the respective target models and the scene. The rejection of a large portion of the candidate models in conjunction with most of the $\bar{R}_{GOF}$ falling on the correct target indicates that the recognition algorithm can readily discriminate the correct target from among the targets in the model library while achieving low false-alarm rates.

In nine out of the twelve scenes, the $\bar{R}_{GOF}$ fell almost entirely on the correct target at $\bar{R}_{GOF}$ levels exceeding 90%. For the remaining three data scenes, the correct target was still assigned the highest $\bar{R}_{GOF}$ value, but a significant portion of the $\bar{R}_{GOF}$ fell on targets other than the target truth. A closer examination of these four scenes reveals that while the $\bar{R}_{GOF}$ did not entirely fall on the correct target, the distribution of $\bar{R}_{GOF}$ values fell almost entirely on a single class of targets that included the target truth.

An example of such a case is the BMP-1 scene that matched the BMP-1 model with an $\bar{R}_{GOF}$ of 0.61 and the BMP-2 model with an $\bar{R}_{GOF}$ of 0.38. Since the BMP-1 and BMP-2 targets have almost identical dimensions and spatial structure, the recognition algorithm was unable to discern the two models from each other. Nonetheless, the scene was recognized to contain a BMP-class vehicle with an $\bar{R}_{GOF}$ of 0.99. Thus we
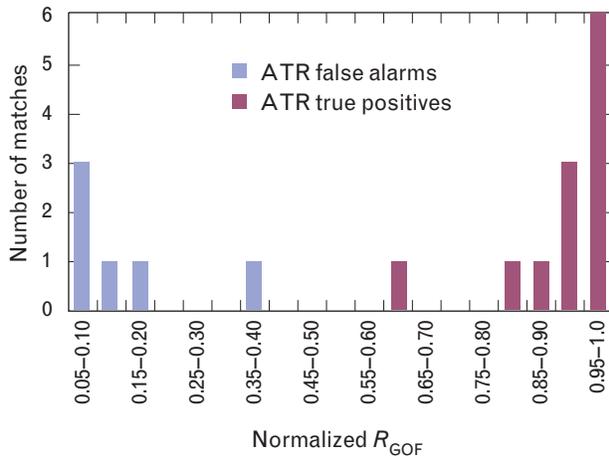
**FIGURE 7.** Distribution of false alarms and true positives over the normalized recognition goodness-of-fit ($R_{GOF}$) value space.

can conclude that the recognition algorithm was able to correctly classify the scene as a BMP with $\bar{R}_{GOF}$ of 0.99 and identify the target as a BMP-1 with a $\bar{R}_{GOF}$ of 0.61.

Another scene that demonstrates correct target classification is the Huntsville T72 scene, where the $\bar{R}_{GOF}$ of the T72 tank model is 0.87 while the $\bar{R}_{GOF}$ of the M1A1 tank model is 0.13. Again, the recognition algorithm correctly classified the scene as a tank with $\bar{R}_{GOF}$ of 1.0 and identified the tank as a T72 with a $\bar{R}_{GOF}$ of 0.87.

Overall, the confusion matrix shows that the recognition algorithm always identified the correct target by assigning the largest $\bar{R}_{GOF}$ value for all twelve recognition tests. To assess recognition performance more clearly, we summarize the data in the confusion matrix into a distribution of false alarms and true positives over the $\bar{R}_{GOF}$ value space, as shown in Figure 7. Given our limited statistics, we have a range of $\bar{R}_{GOF}$ thresholds that allow 100% recognition rate for a 0% false-alarm rate. This range of possible $\bar{R}_{GOF}$ thresholds is determined by the highest $\bar{R}_{GOF}$ false alarm, at 0.38, and lowest $\bar{R}_{GOF}$ true positive, at 0.61. Thus the range of $\bar{R}_{GOF}$ values amounts to a separation of 0.23 in $\bar{R}_{GOF}$ units. This large separation between true positives and false alarms is a good indication of the potential to achieve similar high recognition rates and low false-alarm rates for a larger comparison of scenes.

Table 3 summarizes the average online recognition

timing performance for the twelve scenes. The ATR algorithm was run on a Pentium-4 Xeon 2-GHz machine. In Table 3, the 'average spin-image create time' is the time taken to create the spin images for the automatically selected scene points, the 'average match time' is the average time used to match the scene spin images to each model and generate pose transformations, and the 'average verify time' is the average time taken by the ICP algorithm to verify and refine each scene-to-model comparison. The sum of the stack create time, the average match time, and the average verify time are shown in the column labeled 'total recognition time per model.' The average total time for the twelve scene-to-model library comparisons was approximately two and a half minutes per scene per model.

## Articulated ATR Study

The recognition tests so far have dealt with targets that are represented by solid objects with no articulated components. We now want to extend the ATR algorithm to recognize articulated targets, with multiple movable parts in arbitrary orientations. The main benefit of articulated ATR is that we should have the ability to match an object regardless of the relative position of each of its movable parts (for example, a tank with its turret rotated, or a Scud launcher with its missile at different angular pitches). Furthermore, recognition by parts allows the possibility of recognizing vehicles

**Table 3. ATR Time Performance**

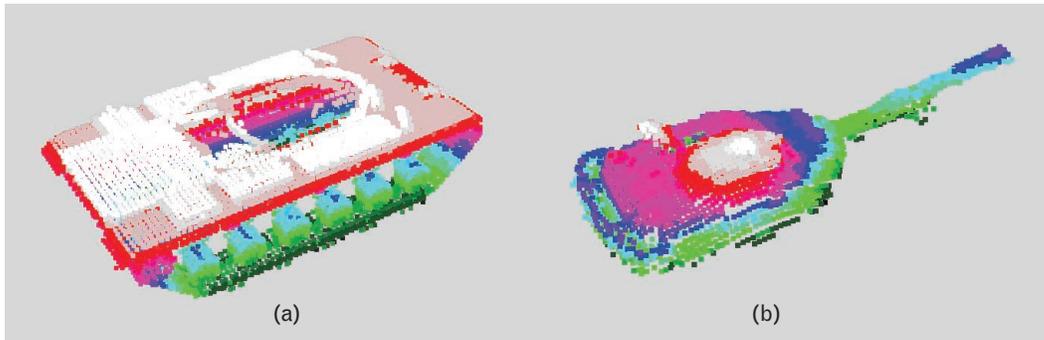| | |
|---|---|
| Average number of scene points | 8570.00 |
| Percentage of scene points selected | 50% |
| Scene resolution (m) | 0.16 |
| Spin-image cylindrical volume (radius, height) | 3,3 |
| Spin-image resolution (pixels × pixels) | 10 × 10 |
| Average spin-image create time (sec) | 14.3 |
| Average match time per model (sec) | 137.50 |
| Average verify time per model (sec) | 3.08 |
| Total recognition time per model (sec) | 142.0 |

**FIGURE 8.** M60 tank parts, color coded by height. (a) The M60 body model; (b) the M60 turret model.

that come in many possible configurations, such as the multipurpose HMMV platform and the myriad of one-of-a-kind technical vehicles encountered in our current military campaigns. Another inherent benefit of articulated ATR is that we can also develop a higher level of tactical awareness by determining the current aim direction of a target's weapon.

We ran a feasibility test to demonstrate articulated ATR on measured Jigsaw data. We created a model library containing two M60 parts—an M60 tank body and an M60 tank turret. Figure 8 shows the two parts in the M60 model library. Figure 9 illustrates the concept of articulated ATR on a scene containing a single-view measurement of an M60 tank with its turret turned by 180°. Figure 10 illustrates a qualitative summary of the results, showing that the correct pose transformation was found for each target part.

To recognize each part in the scene, we consider the measured data present on the other target parts as clutter. For instance, in Figure 10(c) and 10(d), when we are attempting to recognize the M60 turret in the scene,

the measurements on the M60 body act as clutter. Even though the clutter from the M60 body is spatially adjacent to the M60 turret, the recognition algorithm is able to correctly identify the turret and compute a correct pose transformation. The recognition of the body in Figure 10(a) and 10(b) provides another example in which the turret can be considered as close spatial clutter next to the tank body measurement we are attempting to recognize. This successful recognition by parts shows the robustness of the spin-image algorithm to scene clutter, and its potential performance in the development of a fully articulated ATR system.

In the next section we combine our ATR algorithm with an automatic target detection algorithm and show the end-to-end performance of a fully automatic target detection and recognition system.

## Automatic Target Detection in Cluttered Noisy Scenes

Automatic target detection (ATD) was performed by using the general approach of 3D cueing, which deter-
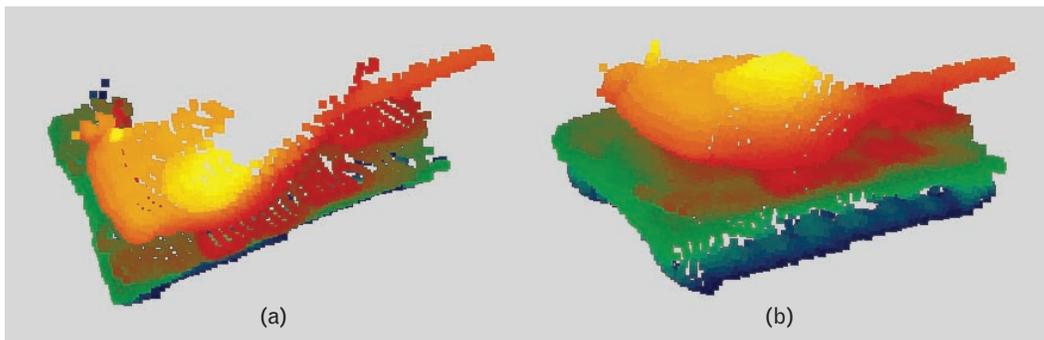


**FIGURE 9.** Single view, color coded by height, of an M60 tank with its turret rotated by 180°. (a) Orthographic view of the scene; ( b) sensor perspective view of the scene.
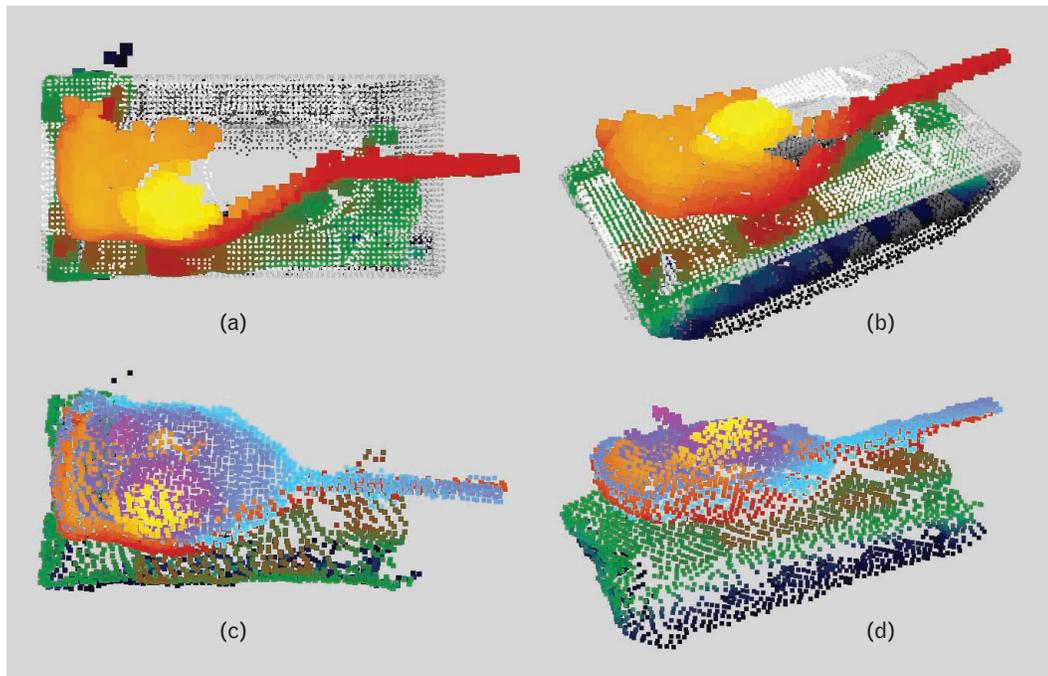
**FIGURE 10.** M60 recognition by parts. (a) Orthographic view of M60 body recognition. The scene points are color coded by height with a green-red-yellow color map, while the M60 model body is color coded by height with a black-to-white color map. (b) Another perspective of the M60 body recognition shows that the correct pose was found in all six degrees of freedom. (c) Orthographic view of the M60 turret recognition. The scene points are again color coded by height with a green-red-yellow color map, while the M60 turret model is color coded by height with a blue-purple color map. (d) Another perspective of the M60 turret recognition shows that the correct pose was found in all six degrees of freedom.

mines and ranks ROIs within a large-scale scene on the basis of the likelihood that these ROIs contain the respective target. Spin-image matching is used to provide a statistical measure of the likelihood that a particular region within the scene contains the target. The detection algorithm is based on the previous work of Carmichael and Hebert et al. [6].

*Detection Algorithm*

The 3D cueing algorithm is tailored for target detection in large-scale terrain scenes. The implemented algorithm can detect and recognize multiple known targets in the scene.

Figure 11 shows a detailed diagram of the automatic target detection and recognition (ATD/R) system for a scene-to-target-model comparison. Following the procedure developed by Carmichael and Hebert et al. [6], we determine ROIs within the scene. The ROI-finding procedure assumes that we test at least one measurement point on the target, although testing as many tar-

get measurements as possible would be optimal. However, in choosing what percentage of points to test from the data set, there is a trade-off between determining the probability to find the target versus the algorithm run-time. On the basis of Carmichael's results and our data, which contains from hundreds to the low thousands of measurements on target, we decided to test between 5% to 10% of the data, with the sampling applied uniformly across the data scene. The ROIs obtained by using the above algorithm can vary drastically in the number of correspondences, correspondence values, and surface area coverage. To discriminate between the various ROIs, we use geometric consistency to remove unlikely correspondences [24]. Each ROI that passes the geometric consistency filter is rated with a detection goodness-of-fit value that corresponds to its likelihood of matching the target of interest. The automatic target detection goodness of fit ($ATD_{GOF}$) value found for ROI $r$ for the comparison of scene $s$ to model $m$ is defined as

**Automatic target detection**

```
Model data ──→ Create model spin images ──── Model spin-image stack ──→ Match model spin images to scene spin image
                                                                         - - - - - - - - -
                                                                         Apply similarity measure filter
Remove ground, trees ──→ Create scene spin images for 10% of points ── Scene spin image ──→ ↑
Extended scene data ──↑                                                  Valid Correspondence?
                        Create scene spin images ←── Neighboring scene points ── Yes
                                                                                 Correspondences ──→ Determine regions of interest
                        Apply geometric consistency filter to each ROI Compute ATD_GOF ←── ROIs
```

**Automatic target recognition**

```
Use ICP to verify and refine transformations ←── ROI transforms ── Group correspondences for each ROI Compute plausible transformations ←── Ordered ROIs
Best refined transform for each ROI ──────────→ Model-to-scene pose and corresponding ATR_GOF value
```
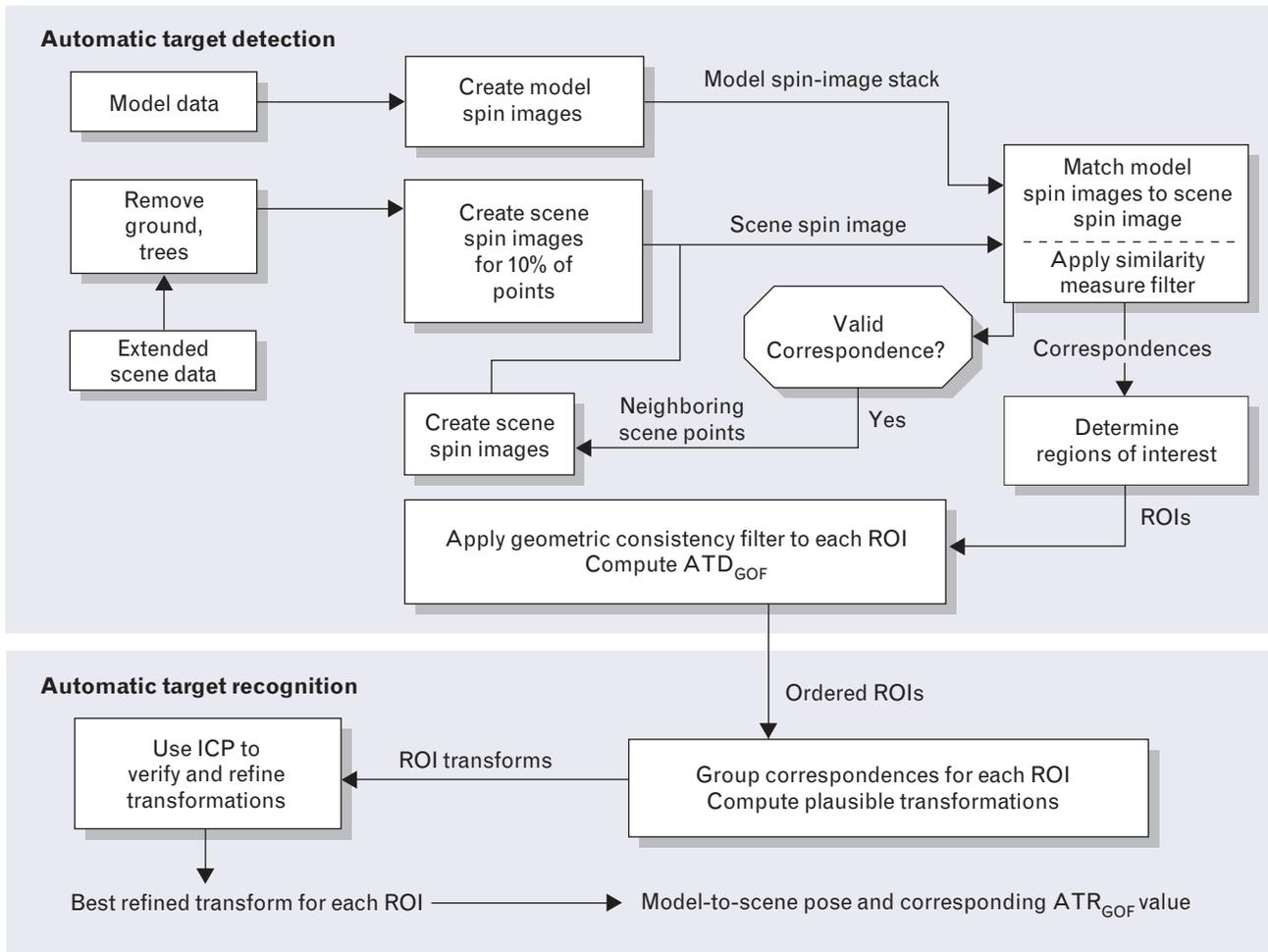
**FIGURE 11.** Process diagram for the automatic target detection (ATD) and automatic target recognition (ATR) system for a scene-to-target-model comparison. The extended scene data are sampled to test at least one measurement on the target (approximately 5% to 10% of points sampled). After matching one or more points on the target with the model, the system explores neighboring points in the scene data and grows a region of interest (ROI). The ROIs are sorted on the basis of their initial likelihood of containing a target, and assigned an ATD goodness-of-fit (ATD$_{GOF}$) value. Each ROI is then sent to the ATR algorithm (previously illustrated in Figure 4), where it is assigned a corresponding ATR goodness-of-fit (ATR$_{GOF}$) value.

$$\mathrm{ATD}_{\mathrm{GOF}}(s,m,r) = \frac{P_r}{Q_r} \sum_{i=1}^{Q_r} C_i \,,$$

where $P_r$ is the number of correspondences in ROI $r$ after the geometric consistency filter, $Q_r$ is the number of correspondences in ROI $r$ before the geometric consistency filter, and $C_i$ is the normalized correlation coefficient value as defined by Johnson et al. [24].

To quantify the detection performance of a scene-to-model library comparison, we normalize the ATD$_{GOF}$ to the maximum ATD$_{GOF}$ value found. The normalized ATD$_{GOF}$ that ROI $r$ in scene $s$ correctly matches

model $i$ in the model library *mlib* is defined as

$$\overline{\mathrm{ATD}}_{\mathrm{GOF}}(s,mlib_i,r) = \frac{\mathrm{ATD}_{\mathrm{GOF}}(s,mlib_i,r)}{\displaystyle\sum_{j=1}^{M}\sum_{k=1}^{N_j} \mathrm{ATD}_{\mathrm{GOF}}(s,mlib_j,k)}\,, \quad (2)$$

where $M$ is the number of models in *mlib*, and $N_j$ is the number of ROIs found for the comparison of scene $s$ to model $mlib_j$. The ROIs are then sorted and queued on the basis of their $\overline{\mathrm{ATD}}_{\mathrm{GOF}}$ value. The recognition algorithm first analyzes the ROI with the best $\overline{\mathrm{ATD}}_{\mathrm{GOF}}$

value before proceeding to the second best ROI, and so on. For each ROI, the recognition algorithm attempts to recognize the model, and then determines a model-to-scene pose and a corresponding $R_{\text{GOF}}$ value (as defined in Equation 1).

To quantify recognition performance of a scene-to-model library comparison, we normalize the $R_{\text{GOF}}$ value to the maximum $R_{\text{GOF}}$ value found. The normalized-to-maximum $R_{\text{GOF}}$ value that ROI $r$ in scene $s$ correctly matches model $i$ in the model library *mlib* is defined as

$$\overline{\text{ATR}}_{\text{GOF}}(s,mlib_i,r) = \frac{R_{\text{GOF}}(r,mlib_i)}{\displaystyle\sum_{j=1}^{M}\sum_{k=1}^{N_j} R_{\text{GOF}}(k,mlib_j)}\;.$$

The end result of the scene-to-model library comparison is a set of ROIs, each matching a target model in a certain pose, along with an $\overline{\text{ATR}}_{\text{GOF}}$ that specifies the level of confidence that the match is correct.

## Results

Five extended terrain scenes recorded with the GEN-III and Jigsaw sensors were used to test the ATD/R system. Each data set contained one or more known targets and covered an area between $25 \times 25$ meters to $100 \times 100$ meters. Target truth in the form of Global Positioning System (GPS) location and target identification was known prior to data collection. Targets in the data set were both out in the open and also underneath heavy canopy cover. Figure 12 shows an orthographic view of the original data sets used for target detection.

Each scene was subsampled by using 20-cm voxels to reduce the computational complexity, and then compared to the target model library. For each ROI
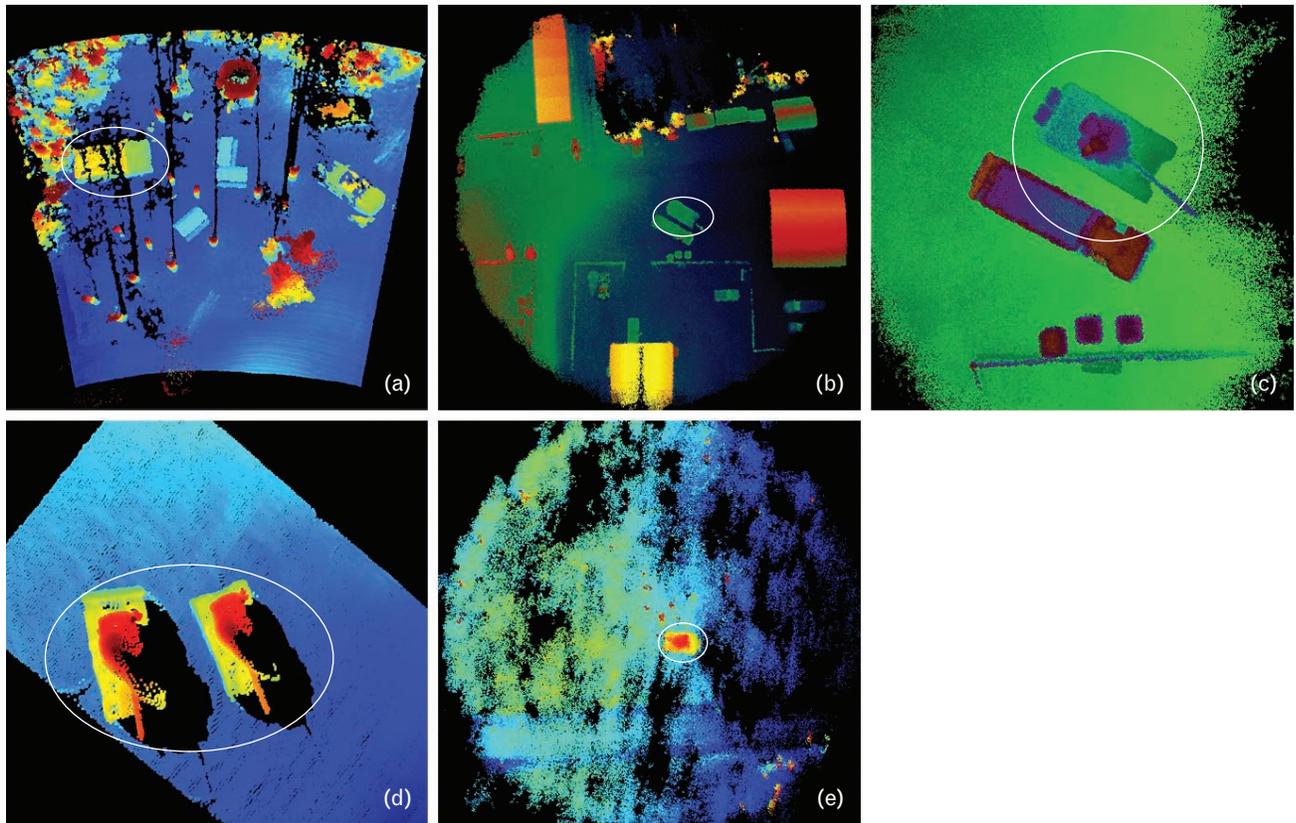


**FIGURE 12.** Orthographic perspective of five large-scale scenes used to test automatic target detection. For some of the data sets, the trees were cropped out to show the obscured target. In each image, the white oval shows the location of the target of interest. (a) GEN-III $25 \times 25$-m measured scene of an HMMV under canopy cover. (b) Jigsaw $100 \times 100$-m measured scene of a T72 in a tank yard from a sensor altitude of 450 m. (c) Jigsaw $25 \times 25$-m measured scene of a T72 in a tank yard from a sensor altitude of 150 m. (d) GEN-III $25 \times 25$-m measured scene of two M60 tanks. (e) Jigsaw $100 \times 100$-m measured scene of a T72 underneath heavy canopy cover, from a sensor altitude of 450 m.
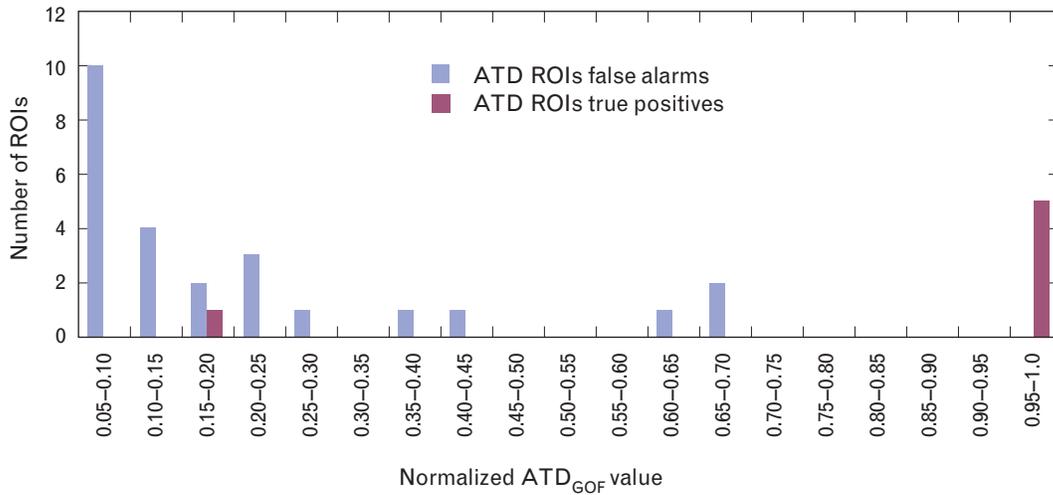
**FIGURE 13.** Distribution of normalized $ATD_{GOF}$ values for the five measured scenes shown in Figure 12. The true positives are shown with magenta color bars, while the false alarms are shown with blue color bars. The ROIs are binned by using a bin size of 0.05 normalized $ATD_{GOF}$ units. For each scene, at least one target instance was detected and mapped to the highest normalized $ATD_{GOF}$ value of 1. In the two-M60 scene in Figure 12(d), the second target instance (which was farther back in the sensor's range) was detected with an $ATD_{GOF}$ of 0.185, which is shown in the 0.15-to-0.20 bin.

found in a scene, we used Equation 2 to compute a value of $\overline{ATD}_{GOF}$. Figure 13 shows the distribution of $\overline{ATD}_{GOF}$ values from all five tested scenes. The distribution is divided between the ROIs that were considered false alarms and the ROIs that were considered true positives. A false alarm is defined as an ROI that matches a target to background clutter or an ROI that incorrectly matches a known scene target to the wrong target model. A true positive is defined as an ROI found for a particular target model that encompasses the measurements of a scene target, and whose target truth matches the respective target model.
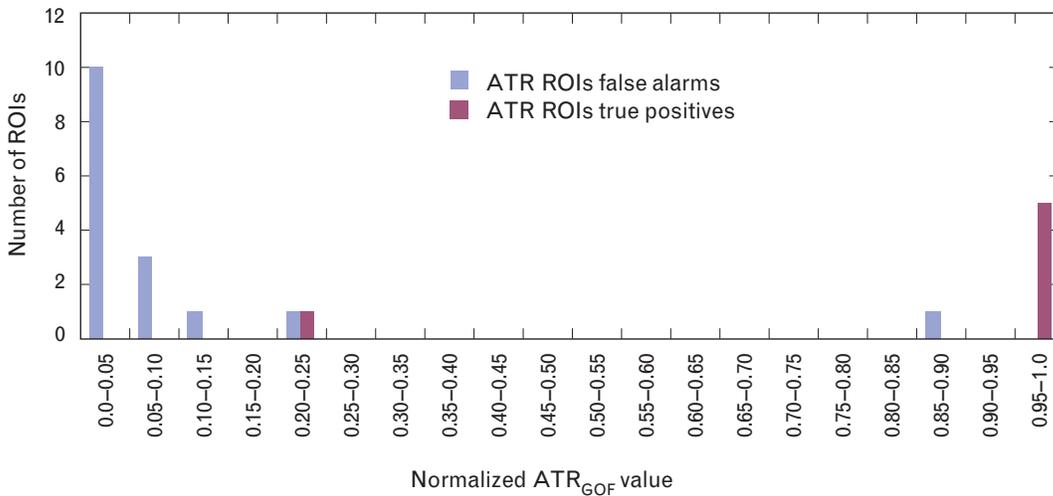


**FIGURE 14.** Distribution of normalized $ATR_{GOF}$ values for the five measured scenes shown in Figure 12. The true positives are shown with magenta color bars, while the false alarms are shown with blue color bars. The ROIs are binned by using a bin size of 0.05 normalized $ATR_{GOF}$ units. For each scene, at least one target instance was detected and mapped to the highest normalized $ATR_{GOF}$ value of 1. In the two-M60 scene in Figure 12(d), the second target instance (which was farther back in the sensor's range) was recognized with an $ATR_{GOF}$ of 0.24, which is shown in the 0.20-to-0.25 bin.

For all five scenes in Figure 12, a true positive ROI had the largest $\overline{\text{ATD}}_{\text{GOF}}$ value, leading to a goodness-of-fit value of one. Thus for all five scenes we were able to correctly detect and identify at least one target instance with a high confidence measure. The M60 scene in Figure 12(d) presented an interesting case, in which two identical M60-type targets existed within the scene. For this single-view scene, the ROI with the highest $\overline{\text{ATD}}_{\text{GOF}}$ of 1.0 fell on the M60 target in the sensor's foreground; the second M60 was also detected, but with a much lower $\overline{\text{ATD}}_{\text{GOF}}$ value of 0.185 (which corresponds in Figure 13 to the true positive under the 0.15–0.20 $\overline{\text{ATD}}_{\text{GOF}}$ bin).

The large difference in $\overline{\text{ATD}}_{\text{GOF}}$ values for the two tanks in the scene is not surprising. The M60 tank in the sensor foreground had about 5318 measurements, while the M60 tank farther down in range from the sensor had about 3676 measurements on its surface. The difference in the number of points is principally due to the difference in angle, resulting in a narrower projected width of the farther M60. Furthermore, the $\overline{\text{ATD}}_{\text{GOF}}$ value is a function of the sum of point-correspondence values and is directly affected by the number of measurements on target. The two-M60 scene presents the following challenge in the detection of multiple instances of a target object within a scene. One of the detected object instances is bound to have a higher signal level than the other objects, which lowers the confidence that the rest of the objects are valid detections of the same target object. In our case, we suspect that the relatively fewer number of data points on the downrange M60 contributes to a normalized $\overline{\text{ATD}}_{\text{GOF}}$ confidence value that is smaller than the $\overline{\text{ATD}}_{\text{GOF}}$ value of the foreground M60.

If we ignore the low-$\overline{\text{ATD}}_{\text{GOF}}$ true-positive result from the M60 scene, Figure 13 shows a good separation between the distributions of false alarms and true positives. The two distributions have a separation of about 0.33 in the $\overline{\text{ATD}}_{\text{GOF}}$ value space. This difference indicates that we can always detect and identify the correct target from the library of known targets. With a separation of 0.33 in the $\overline{\text{ATD}}_{\text{GOF}}$ value space, a detection threshold can readily be set between the highest false alarm (at 0.671) and the lowest of the remaining true positives (at 1.0). Thus, even as a stand-alone algorithm, the ATD system works exceptionally well.

## Combining ATD and ATR

We now show the results of ATD coupled with ATR. Figure 14 shows the distribution of $\overline{\text{ATR}}_{\text{GOF}}$ values obtained after we ran the ATR algorithm on the detected ROIs. From the distributions, we can discern that most of the true positives are mapped to the highest $\overline{\text{ATR}}_{\text{GOF}}$ value of one. Again, the multiple M60 targets presented a challenge with the background M60 tank mapping to a normalized $\overline{\text{ATR}}_{\text{GOF}}$ of 0.24, slightly higher than the 0.185 $\overline{\text{ATD}}_{\text{GOF}}$ value for the background tank. There is also a significant improvement in the distribution of false alarms and true positives in the $\overline{\text{ATR}}_{\text{GOF}}$ value space as compared to the $\overline{\text{ATD}}_{\text{GOF}}$ value space. Most of the ATD false alarms have been remapped from an $\overline{\text{ATD}}_{\text{GOF}}$ range of 0 to 0.67 to an $\overline{\text{ATR}}_{\text{GOF}}$ range of 0 to 0.24. The remapping of false alarms from higher $\overline{\text{ATD}}_{\text{GOF}}$ values to lower $\overline{\text{ATR}}_{\text{GOF}}$ values further increases the separation between the distribution of false alarms and true positives. The larger separation between the majority of false alarms and true positives represents an improvement in our ability to discern the correct target from background clutter and other known targets. Therefore, the $\overline{\text{ATR}}_{\text{GOF}}$ value space is an improvement over the $\overline{\text{ATD}}_{\text{GOF}}$ value space.

Table 4 shows the time performance of the entire ATD and ATR system. The ATD/R system was run on an Intel Pentium-4 Xeon at 2 GHz. In the table, 'stack create time' is the time taken to create the spin-image stack of the scene. The 'average ATD+ATR time per model' is the time used to detect ROIs for a model, and recognize whether the ROI is a valid target model instance. The 'average ATD+ATR time per model' also includes the contribution of the time taken to create the scene spin-image stack, weighted down by the number of models in the library, since the scene stack is computed only once and used for all the following target model comparisons. Overall, we achieved a recognition time of approximately one and a half minutes per model.

In summary, our new ATD/R algorithm has demonstrated very good detection and identification accuracy, as well as time performance. Given its timing and accuracy performance, this ATD/R system may have significant practical value to a human operator for aided target recognition under battlefield conditions.

**Table 4. ATD and ATR System Time Performance**

| Scene description | Total number of scene points before ground removal | Total number of scene points after ground removal | Target percentage of scene by volume, based on original scene with ground and canopy | Percent of scene points selected to correlate to models | Scene resolution (m) | Spin-image size (number of pixels) | Spin-image stack create (sec) | Average ATD + ATR time per model, based on the ten-model library * (sec) |
|---|---|---|---|---|---|---|---|---|
| **HMMV scene** Huntsville June 2003 C8-F1-P10 | 192,097 | 26,318 | 0.76 | 100% | 0.25 | 25 (5×5) | 59.8 | 120.72 |
| **M60s scene** Eglin 2001 | 48,997 | 8995 | 9.18 | 100% | 0.16 | **100** (10×10) | 15.6 | 497.86 |
| **Tank yard** (450-m altitude) Huntsville Dec 2002 C20-F02-P05 | 575,938 | 35,157 | 0.59 | 100% | 0.26 | 25 (5×5) | 40.9 | 219.26 |
| **T72 under canopy** (450-m altitude) Dec 2002 Huntsville C20-F02-P07 | 312,189 | 10,293 | 2.41 | 100% | 0.18 | 25 (5×5) | 29.9 | 45.33 |
| **Tank yard** (150-m altitude) Huntsville Dec 2002 C20-F01-P3 | 32,750 | 7286 | 8.99 | 100% | 0.24 | 25 (5×5) | 7.62 | 30.69 |

* Average ATD+ATR Time for 20 cm subsampled scenes (seconds) = 104.00

## Conclusions

In this research, we developed and implemented a fully automated target detection and recognition system that uses geometric shape and size signatures from target models to detect and recognize targets under heavy canopy and camouflage cover in extended terrain scenes.

The ATD/R system performance was demonstrated on five measured scenes with targets both out in the open and under heavy canopy cover, where the target occupied between 1% to 10% of the scene by volume. The ATR section of the system was successfully demonstrated for twelve measured data scenes with targets both out in the open and under heavy canopy and camouflage cover. Correct target identification was also demonstrated for targets with multiple movable parts that are in arbitrary orientations. We achieved a high recognition rate (over 99%) along with a low false-alarm rate (less than 0.01%).

The major contribution of this research is that we proved that spin-image–based detection and recognition is feasible for terrain data collected in the field with a sensor that can be used in a tactical situation. We also demonstrated recognition of articulated objects, with multiple movable parts. Considering the detection and recognition performance, the ATD/R system can have

significant practical value to a human operator for aided target recognition under battlefield conditions.

Immediate benefits of the presented work will be in the area of automatic target recognition of military ground vehicles, where the vehicles of interest may include articulated components with variable position relative to the body, and come in many possible configurations. Other application areas include human detection and recognition for homeland security.

# REFERENCES

1. R.M. Marino, T. Stephens, R.E. Hatch, J.L. McLaughlin, J.G. Mooney, M.E. O'Brien, G.S. Rowe, J.S. Adams, L. Skelly, R.C. Knowlton, S.E. Forman, and W.R. Davis, "A Compact 3D Imaging Laser Radar System Using Geiger-Mode APD Arrays: System and Measurements," *SPIE* **5086**, 2003, pp. 1–15.

2. R.M. Heinrichs, B.F. Aull, R.M. Marino, D.G. Fouche, A.K. McIntosh, J.J. Zayhowski, T. Stephens, M.E. O'Brien, and M.A. Albota, "Three-Dimensional Laser Radar with APD Arrays," *SPIE* **4377**, 2001, pp. 106–117.

3. J.A. Ratches, C.P. Walters, R.G. Buser, and B.D. Guenther, "Aided and Automatic Target Recognition Based upon Sensory Inputs from Image Forming Systems," *IEEE Trans. Patt. Anal. Mach. Intell.* **19** (9), 1997, pp. 1004–1019.

4. M. Wellfare and K. Norris-Zachery, "Characterization of Articulated Vehicles Using Ladar Seekers," *SPIE* **3065,** 1997, pp. 244–254.

5. J.-Y. Dufour and V. Martin, "Active/Passive Cooperative Image Segmentation for Automatic Target Recognition," *SPIE* **2298,** 1994, pp. 552–560.

6. O. Carmichael and M. Hebert, "3D Cueing: A Data Filter for Object Recognition," *IEEE Conf. on Robotics and Automation* **2,** Detroit, 10–15 May 1999, pp. 944–950.

7. G.D. Arnold, K. Sturtz, and I. Weiss, "Detection and Recognition in LADAR Using Invariants and Covariants," *SPIE* **4379,** 2001, pp. 25–34.

8. Y.-T. Zhou and D. Sapounas, "An IR/LADAR Automatic Object Recognition System," *SPIE* **3069,** 1997, pp. 119–128.

9. S. Grossberg and L. Wyse, "A Neural Network Architecture for Figure-Ground Separation of Connected Scenic Figures," *Neural Netw.* **4** (6), 1991, pp. 723–742.

10. Z. Ying and D. Castanon, "Statistical Model for Occluded Object Recognition," *Proc. 1999 Int. Conf. on Information Intelligence and Systems, Bethesda, Md., 31 Oct.–3 Nov. 1999,* pp. 324–327.

11. F. Sadjadi, "Application of Genetic Algorithm for Automatic Recognition of Partially Occluded Objects," *SPIE* **2234,** 1994, pp. 428–434.

12. M.A. Khabou, P.D. Gader, and J.M. Keller, "LADAR Target Detection Using Morphological Shared-Weight Neural Networks," *Mach. Vis. Appl.* **11** (6), 2000, pp. 300–305.

13. M. Hebert and J. Ponce, "A New Method for Segmenting 3-D Scenes into Primitives," *Proc. 6th Int. Conf. on Pattern Recognition* **2,** *Munich, 19–22 1982,* pp. 836–838.

14. R. Hoffman and A.K. Jain, "Segmentation and Classification of Range Images," *IEEE Trans. Patt. Anal. Mach. Intell.* **9** (5), 1987, pp. 608–620.

15. D.L. Milgram and C.M. Bjorklund, "Range Image Processing: Planar Surface Extraction," *Proc. 5th Int. Conf. on Pattern Recognition* **1,** *Miami Beach, Fla., 1–4 Dec. 1980,* pp. 912–919.

16. J. Huang and C.-H. Menq, "Automatic Data Segmentation for Geometric Feature Extraction from Unorganized 3-D Coordinate Points," *IEEE Trans. Robot. Autom.* **17** (3), 2001, pp. 268–279.

17. I.K. Park, I.D. Yun, and S.U. Lee, "Automatic 3-D Model Synthesis from Measured Range Data," *IEEE Trans. Circuits Syst. Video Technol.* **10** (2), 2000, pp. 293–301.

18. D. Cobzas and H. Zhang, "Planar Patch Extraction with Noisy Depth Data," *Third Int. Conf. on 3-D Digital Imaging and Modeling, Quebec City, Canada, 28 May–1 June 2001,* pp. 240–245.

19. F. Stein and G. Medioni, "Structural Indexing: Efficient 3-D Object Recognition," *IEEE Trans. Patt. Anal. Mach. Intell.* **14** (2), 1992, pp. 125–145.

20. O. Carmichael, D.F. Huber, and M. Hebert, "Large Data Sets and Confusing Scenes in 3-D Surface Matching and Recognition," *Proc. Second Int. Conf. on 3-D Digital Imaging and Modeling, Ottawa, 4–8 Oct. 1999,* pp. 358–367.

21. V. Shantaram and M. Hanmandlu, "Contour Based Matching Technique for 3D Object Recognition," *Proc. Int. Conf. on Information Technology: Coding and Computing* **3,** *Las Vegas, Nev., 8–10 Apr., 2002,* pp. 274–279, 2002.

22. C. Dorai and A.K. Jain, "Shape Spectra Based View Grouping for Free-Form Objects," *Proc. Int. Conf. on Image Processing* **3,** *Washington, D.C., 22–26 Oct. 1995,* pp. 340–343.

23. S. Yamany and A. Farag, "3D Objects Coding and Recognition Using Surface Signatures," *Proc. 15th Int. Conf. on Pattern Recognition* **4,** *Barcelona, 3–7 Sept. 2000,* pp. 571–574.

24. A. Johnson, "Spin-Images: A Representation for 3-D Surface Matching," Ph.D. thesis (Robotics Institute, Carnegie Mellon University, Pittsburgh, 1997).

25. A.N. Vasile, "Pose Independent Target Recognition System Using Pulsed Ladar Imagery," Master of Engineering thesis (Electrical Engineering and Computer Science, MIT, Cambridge, Mass., 2003).

26. J. Schroeder, "Extended Abstract on Automatic Target Detection and Recognition Using Synthetic Aperture Radar Imagery," Cooperative Research Centre for Sensor Signal and Information processing (CSSIP) SPRI building, Mawson Lakes Boulevard Mawson Lakes, South Australia, <http://www.ips.gov.au/IPSHosted/NCRS/wars/wars2002/proceedings/invited/print/schroeder.pdf>.

**ALEXANDRU N. VASILE**
is an associate staff member in the Active Optical Systems group. He received S.B. and M.Eng. degrees in electrical engineering and computer science from MIT. His research interests include artificial intelligence, computer vision, biomedical imaging, and imaging algorithm development. He joined Lincoln Laboratory in 2000 as an undergraduate student, with the ultimate goal of doing a Master's thesis in computer vision with an emphasis on 3D target visualization and automatic target recognition. As an undergraduate student he also helped develop an image query system, based on color harmony, to improve the efficiency of image search methods. He has worked at the MIT Media Laboratory; Radiation Monitoring Devices, Inc.; and the Electro-Optics Technology Center at Tufts University.

**RICHARD M. MARINO**
is a senior staff member in the Active Optical Systems group. He received a B.S. degree in physics from Cleveland State University, and an M.S. degree in physics and a Ph.D. degree in high-energy physics from Case Western Reserve University. He joined Lincoln Laboratory as a staff member in the Laser Radar Measurements group, and later joined the Systems and Analysis group. One of his most significant achievements has been his pioneering leadership in the development of a 3D imaging laser radar with photon counting sensitivity. He has also worked at the Millimeter Wave Radar (MMW) and the ARPA-Lincoln C-band Observables Radar at the Kwajalein Missile Range in the Marshall Islands. While there, he was a mission test director at MMW and worked on range modernization plans. In 1997 he joined the Sensor Technology and Systems group of the Aerospace division and relocated its Silver Spring, Maryland, location to join the National Polar-Orbiting Operational Environmental Satellite System (NPOESS)/Integrated Program Office (IPO). At the IPO, he was lead technical advisor for the NPOESS Cross-Track Infrared Atmospheric Sounder Instrument (CrIs). He returned to Lincoln Laboratory in Lexington in 1999 and is again working on the development of 3D imaging laser-radar systems.